# Seaview Survey Photo-quadrat and Image Classification Dataset

González-Rivero, Manuel, Rodriguez-Ramirez, Alberto, Beijbom, Oscar, Ganase, Anjani, Dalton, Peter, Kennedy, Emma V., Neal, Benjamin P., Vercelloni, Julie, Bongaerts, Pim, Bryant, Dominic E.P., Brown, Kristen, Kim, Catherine, Radice, Veronica Z., Lopez-Marcano, Sebastian, Dove, Sophie, Bailhache, Christophe, Beyer, Hawthorne L. & Hoegh-Guldberg, Ove

Global Change Institute & School of Biological Sciences, The University of Queensland, Australia

**Project**: The XL Catlin Seaview Survey Project was developed as a collaboration between The University of Queensland and ocean conservation non-profit Underwater Earth.

**Project description**: The goal of the XL Catlin Seaview Survey Project was to conduct rapid, detailed yet globally distributed scientific surveys of coral reefs to support research and conservation. The focus was on surveying shallow reefs, typically located at 10 m depth, around the world: 860 surveys across the Atlantic (mainly Caribbean), the Indian Ocean, the Pacific Ocean, the Great Barrier Reef, and South-East Asia were surveyed from 2012-2018. A custom camera and sensor system (the "SVII") was developed by Underwater Earth and adapted for science by The University of Queensland in order to collect high quality photographs of the reef at regular intervals along sampling transects that are usually 1.5-2.0 km in length. The system also collected camera altitude from the reef substrate and GPS data to assist in the processing of the photographs and assessment of coral reefs.



*Figure 1. Dr. Manuel González-Rivero piloting the "SVII" camera system. © Underwater Earth / XL Catlin Seaview Survey / Christophe Bailhache*

**Dataset description:** The primary scientific datasets arising from the project is the "Seaview Survey Photo-quadrat and Image Classification Dataset", consisting of: (1) over one million standardised, downward-facing "photo-quadrat" images covering approximately 1m$^2$ of the sea floor; (2) human-classified annotations that can be used to train and validate image classifiers; and (3) benthic cover data arising from the application of machine learning classifiers to the photo-quadrats. Photo-quadrats were collected between 2012 and 2018 at 860 transect locations around the world, including: the Caribbean and Bermuda, the Indian Ocean (Maldives, Chagos Archipelago), the Coral Triangle (Indonesia, Philippines, Timor-Leste, Solomon Islands), the Great Barrier Reef, Taiwan and Hawaii. Multi-temporal images exist for some sites in the Great Barrier Reef, Indonesia, Maldives, Philippines, Hawaii and Taiwan.



*Figure 2. Locations of the 860 Seaview Survey transects (red dots). Background imagery reproduced from the GEBCO World Map (www.gebco.net)*

## Methodology

Photographs were taken using Canon 5D MII cameras without artificial light and using a fisheye lens to maximise light capture. Lens distortion and colour balancing was initially applied to RAW format images using Photoshop (Adobe Systems, California, USA), a process that was later auto-mated using Photoshop and ImageMagick (www.imagemagick.org). Colour balancing images improved the performance of image classifiers, but different approaches to lens correction and colour balancing were found to have no effect on their performance. Photo-quadrats were extracted from the central area of these corrected images using the altitude data logged by an acoustic depth sounder to standardise photo-quadrat area to approximately 1m$^2$. From one single downward facing image collected, several photo-quadrats can be extracted depending on the distance between the camera and the reef surface. Altitude estimates, combined with the attributes of the camera and lens (such as sensor dimensions, focal length, and geometric distortion of the lens) was used to estimate the angular field of view of the camera and, with this, the footprint of each image. Considering a margin around the image to eliminate the vignetting effect from post-processing, the size of each image (i.e., footprint) determined the number of photo-quadrats (1m$^2$ cropped sections) to extract from the original image. Photo-quadrats were only generated from images within a range of altitude between 0.5 and 2 m to ensure a consistent spatial resolution for each image (approx. 10 pixels/cm).
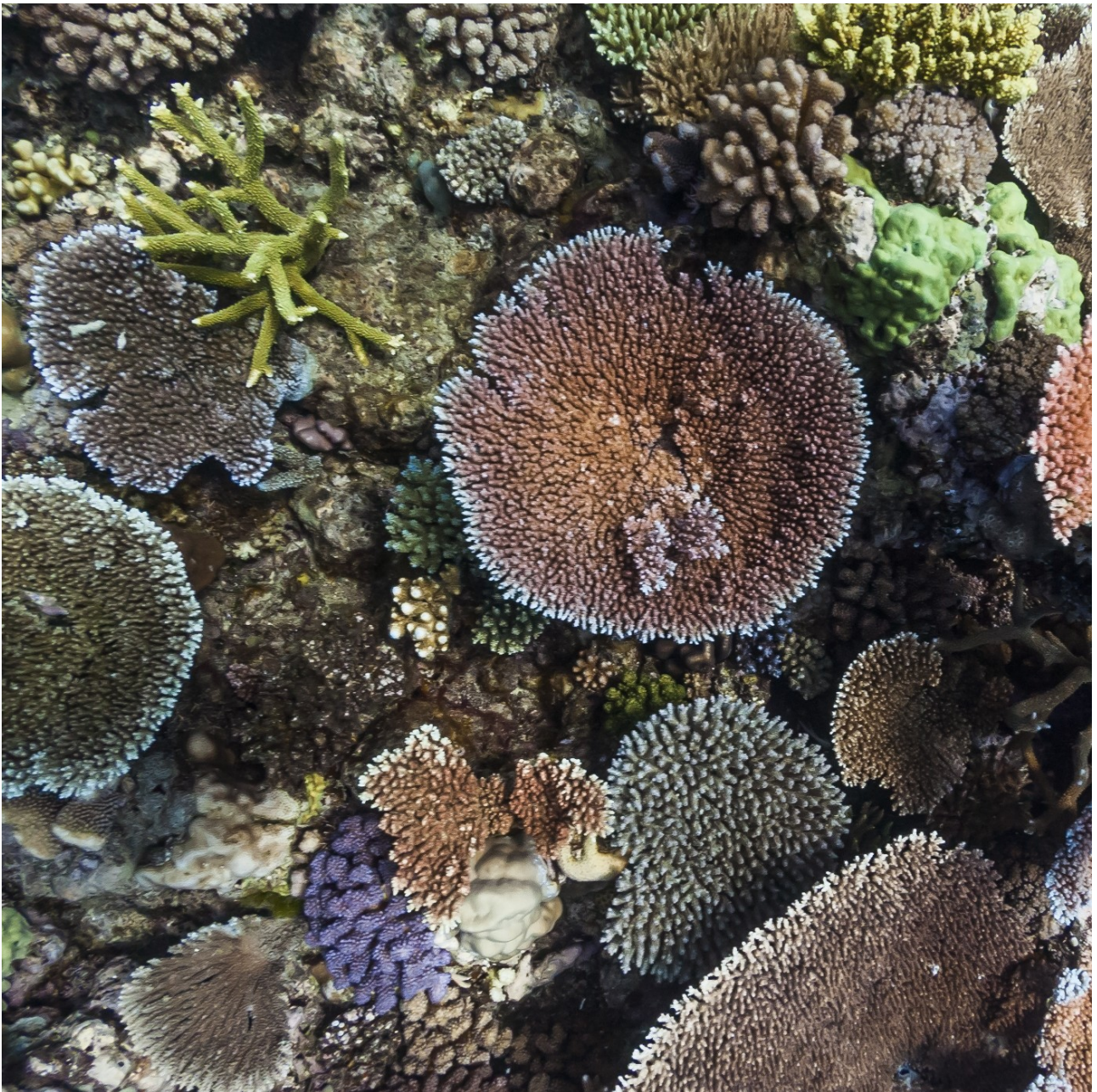
*Figure 3. An example of a photo-quadrat image. This is representative of the better quality images in the data collection. Factors such as turbidity, light and wind conditions, the height of the camera above the reef, and the speed of the camera can degrade image quality.*
© The University of Queensland / Underwater Earth / XL Catlin Seaview Survey

Deep learning algorithms (specifically a VGG-D 16 network architecture) were used to estimate the benthic cover from each photo-quadrat, using random point annotations. This method, also referred as random point count, is commonly used in many population estimation applications using photographic records. In random point annotations, the relative cover or abundance of each class is defined by the number of points classified as such relative to the total number of observed points on the image. In order to achieve automated random point annotation, we converted each image to a set of patches cropped out around each given point location. The patch area to crop around each point was set fix to 224 x 224 to align with the pre-defined image input size of the VGG architecture. Each cropped patch was then classified independently and then the relative abundance for each of the benthic classifications was the ratio between the number of patches classified for a given class by the total number of patches evaluated in an image.

In this study, a VGG network pre-trained on ImageNet was initialised and then fine-tuned using the random selection of images manually annotated as the training dataset. This fine-tuning exercise ran for 40 K iterations, where the final classification was contrasted against a validation set of images, an independent 20% subset from the original set of training images (more details in https://github.com/mgonzalezrivero/reef_learning). To estimate the benthic cover estimation errors from automated classification, an independent set of images within transects where manually annotated and contrasted against the automated cover estimations for the same region. A different deep learning network was trained for each country, with the exception of Easter Atlantic where all the images per country were used to train a single network.
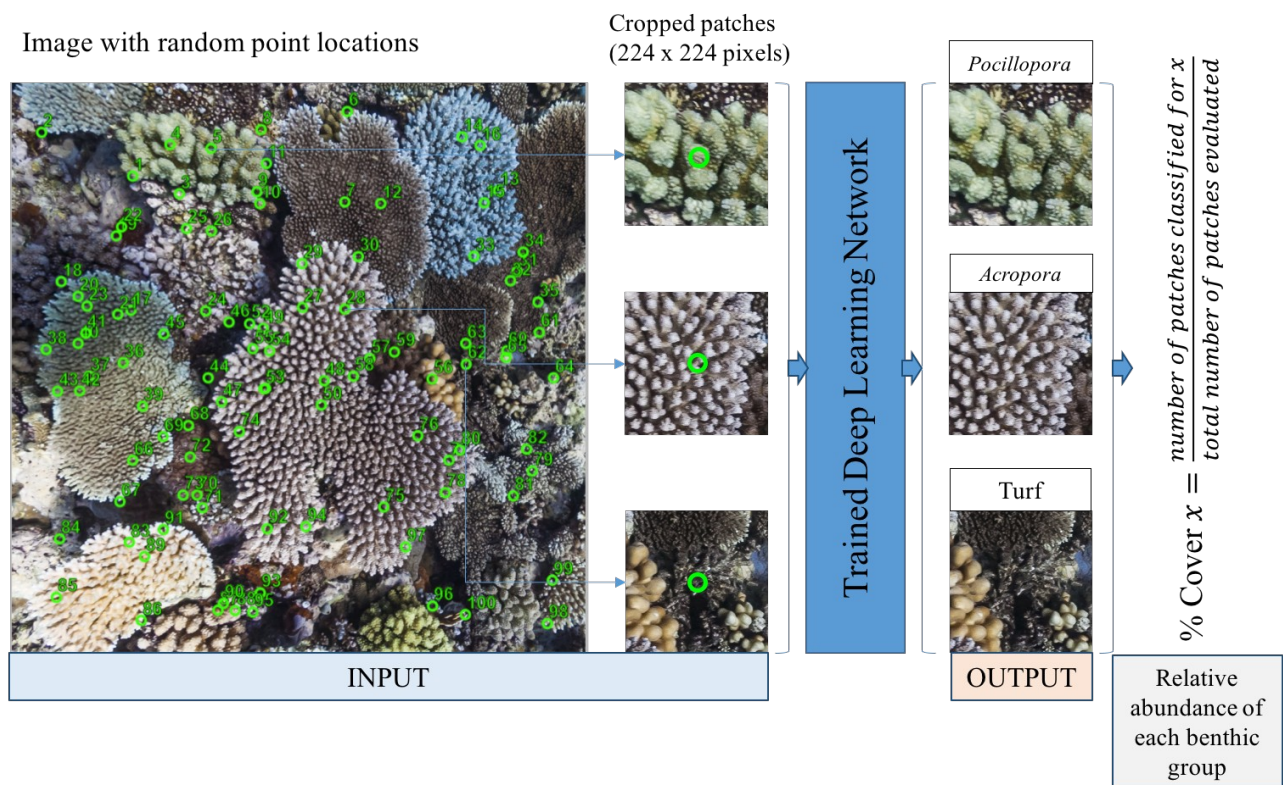


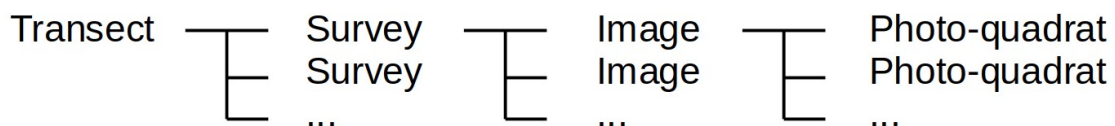*Figure 4. Illustration of the image training and classification process.*



*Figure 5. Hierarchical structure of photo-quadrat dataset. A transect refers to a location that can be surveyed multiple times (e.g. in different years). Each survey involves the collection of many images. From each image one or more standardised photo-quadrat images can be extracted.*

## Dataset organisation

### Image Data

**Organisation and size:** The photo-quadrat images are in JPEG format, with files grouped by survey. There are 1.1 million photo-quadrat images among 860 surveys, totalling 1.5 TB. Each survey contains a median of 1075 quadrat images (range: 94-5754), with a median transect folder size of 1.7 GB (range: 0.1-6.2 GB). To facilitate data curation the photo-quadrats associated with each survey are contained in 860 separate archive files (.zip files) that range from 0.1-6.2 GB GB in size and are located within folder "photo-quadrats". See Figure 5 for an explanation of the hierarchical structure of the dataset.

**Survey naming convention:** Survey names consist of the combination of four information components separated by underscores: the ocean in which the transect occurs (ATL=Atlantic Ocean, IND=Indian Ocean, PAC=Pacific Ocean), the three letter country code of the Exclusive Economic Zone within which the transect occurs (e.g. AUS=Austalia), the unique five digit survey ID number, and the year and month of the survey formatted as YYYYMM. Some examples of this naming convention are: ATL_ABW_17006_201304, ATL_BLZ_20047_201307, IND_CHA_36001_201502, PAC_AUS_47035_201710, PAC_IDN_64047_201806, PAC_USA_44036_201606. Each of these survey folders has been compressed into a zip file (e.g. ATL_ABW_17006_201304.zip).

**Survey preview images:** Low resolution previews of all of the photo-quadrats in each survey (tiled into one JPEG image per survey – see Figure 6) have been zipped into archives grouped by ocean and country (e.g. PAC_AUS referes to the Pacific Ocean and Australia). These files are available in the "survey-previews" folder. These preview images provide a coarse overview of the content of the photo-quadrats that may help users to identify which surveys they are most interested in downloading.
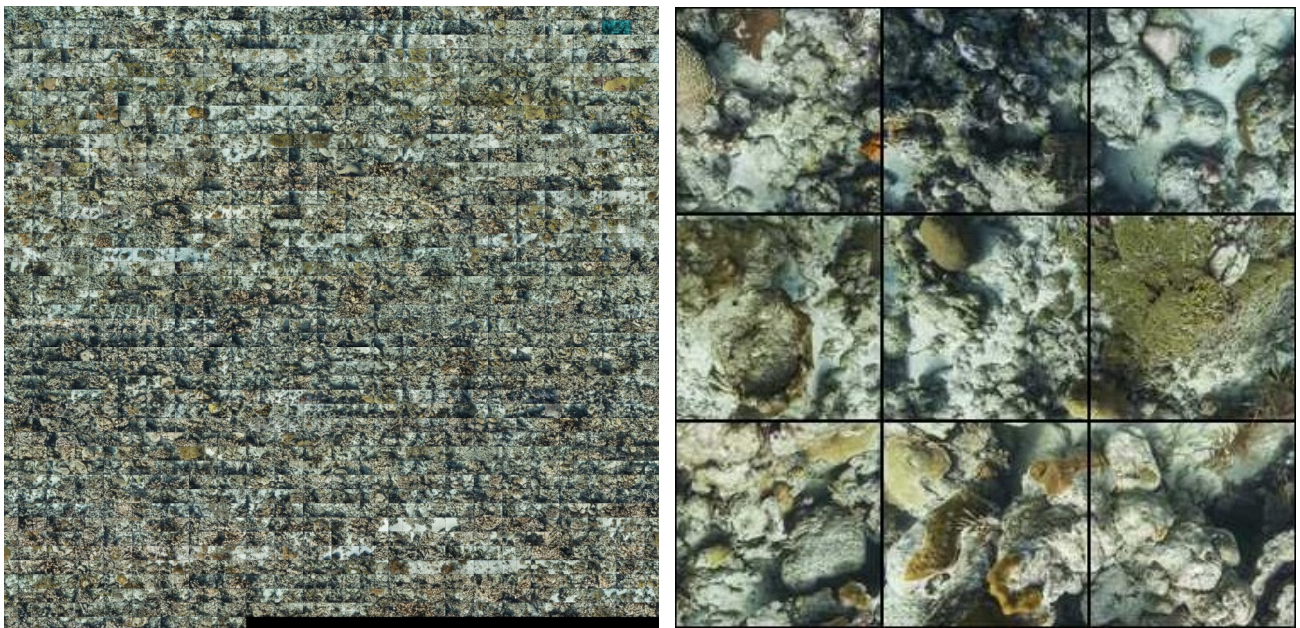


*Figure 6. Low resolution preview images, arranged into a single montage image (left) of the photo-quadrats in each survey are included to provide an indication of the content of the pictures in the survey (right: a blow-up of a 3 x 3 block of photo-quadrat images from the montage image on the left).*

**Annotated images:** Coral reef images used in the development of the training and validation data (n = 11,387 images) are grouped according to classification schema as different schema were used in each major region. These images are located in the "annotated-images" folder, which contains the following region folders:

| Folder | Region | Folder size | Number of images |
|---|---|---|---|
| ATL | Atlantic | 2.1 GB | 1407 |
| IND_CHA | Indian Ocean, Chagos Archipelago | 0.9 GB | 686 |
| IND_MDV | India Ocean, Maldives | 2.5 GB | 1612 |
| PAC_AUS | Pacific Ocean, Australia | 2.7 GB | 2657 |
| PAC_USA | Pacific Ocean, USA (Hawaii) | 2.9 GB | 1153 |
| PAC_IDN_PHL | Pacific Ocean, Indonesia and Philippines | 2.7 GB | 1638 |
| PAC_SLB | Pacific Ocean, Solomon Islands | 0.8 GB | 732 |
| PAC_TWN | Pacific Ocean, Taiwan | 0.7 GB | 638 |
| PAC_TLS | Pacific Ocean, Timor-Leste | 1.4 GB | 864 |

*Table 1. Organisation of the images associated with the annotations data, which is used to train and validation image classification algorithms. Label sets differ among the regions, hence the stratification of the annotation images and tabular data into separate folders.*


## Tabular data

The tabular data is organised into a set of related and hierarchical tables (compressed, comma-delimited CSV files located within an "survey_image_data" folder):

**1. Survey table** (860 rows, seaviewsurvey_surveys.csv). A survey refers to a transect conducted at a specific location and date, hence a single transect can be surveyed repeatedly over time. The table has the following fields:

- surveyid: A unique survey ID code representing data collection at one location (a transect location) at one point in time. The survey IDs are unique in this table, and these IDs are used in the folder naming convention. Note that there may be multiple survey IDs associated with a transect ID in the subset of cases in which multi-temporal surveys were conducted.
- transectid: The five digit transect identification code. A transect ID will appear more than once in this field if the transect has been sampled on more than one occasion.
- surveydate: the date (YYYYMMDD) on which the survey was completed
- ocean: the three letter code representing the ocean within which the survey occurred (ATL=Atlantic Ocean, IND=Indian Ocean, PAC=Pacific Ocean)
- country: the three letter country code of the Exclusive Economic Zone within which the survey occurs (e.g. AUS=Austalia)
- foldername: the full name of the folder associated with the survey (e.g. " PAC_AUS_47035_201710")
- lat_start, lng_start: the latitude and longitude of the start of the survey
- lat_end, lng_end: the latitude and longitude of the end of the survey

- pr_hard_coral, pr_algae, pr_soft_coral, pr_oth_invert, pr_other: The proportional cover of the highest-level of the reef cover classification categories, averaged across all photo-quadrats in the survey

**2. Quadrats table** (~1.1 million rows, seaviewsurvey_quadrats.csv). Photographs are taken approximately once every two seconds along a transect.

- surveyid: The five digit survey ID number that is in the survey table. There is a one:many relationship between the survey table and the quadrat table.
- imageid: The image ID number representing the image from which this quadrat was extracted. Note that there can be duplicates in this field because more than one quadrat can be derived from some images.
- quadratid: The unique quadrat ID number that is also the filename (".jpg" must be added to the quadrat ID to derive the filename with the extension).

**3. Annotations tables** (~55.2 million rows. seaviewsurvey_annotations.csv). An annotation refers to the classification of a pixel in a quadrat image into one of the label categories for that region (algae, various coral groups or species, sand, etc). This classification is performed by an expert human, and this dataset is intended to be used for image classification and training purposes.

- quadrat ID: The quadrat ID number. There is a one:many relationship between the quadrat table and the annotations table because each quadrat can have many annotations.
- x, y: the coordinates of the pixel that has been classified, in reference to the top left of the image with an origin 1,1 index
- label: the short code representing the features of interest

Different label sets were used in different regions, reflecting differences in species composition and the ability to distinguish coral cover types. The annotation tables for each region are named "annotations_" with the region code (see Table 1) appended. Specifically:

annotations_ATL.csv
annotations_IND_CHA.csv
annotations_IND_MDV.csv
annotations_PAC_AUS.csv
annotations_PAC_IDN_PHL.csv
annotations_PAC_SLB.csv
annotations_PAC_TLS.csv
annotations_PAC_TWN.csv
annotations_PAC_USA.csv

Note that the filename of the image associated with each annotation record is the quadrat ID number with ".jpg" appended to the end, and the subfolder in the annotations folder containing the image can be inferred from the filename of the annotations file. For example, all the images referred to in the annotations_PAC_AUS.csv file are in the annotations/PAC_AUS folder.

**4. Reef Cover tables by region** (~1.1 million rows across all tables, seaviewsurvey_reefcover_[region].csv). The reef cover tables contain the summary classification data for each of the photo-quadrats, represented as the estimated percent cover of each of the cover types. The reef cover data is segregated by region because the classification schemes differ among these regions (Caribbean, the Coral Triangle and Indo-Pacific, the Great Barrier Reef, Hawaii, and the Indian Ocean). The fields include:

- surveyid: The five digit transect ID number that is in the survey table. There is a one:many relationship between the surveys table and the quadrat table.
- imageid: The image ID number representing the image from which this quadrat was extracted (".jpg" must be added to the image ID to derive the filename with the extension).
- quadrat ID: The quadrat ID number. There is a one:many relationship between the quadrat table and the annotations table because each quadrat can have many annotations. ".jpg" must be added to the quadrat ID to derive the filename with the extension.
- latitude and longitude of quadrat
- [classname]: The proportional cover of the specified class name within the photo-quadrat

**5. Label key table** (228 rows, seaviewsurvey_labelsets.csv). This table provides a more detailed explanation of what the annotation labels mean.

- region: the region identifier for the label set
- label: the short code representing the features of interest
- func_group: a more general categorisation of the features
- label_name: the full name of the label
- merged_label: alternative short code in order reduce the complexity of the labelset. Some related labels were merged
- merged_name: the full name of the merged label
- description_examples: a brief description of what the labels represent, with examples

## Use case: Accessing reef cover data

Users primarily interested in accessing percent cover data for each of the 860 surveys need to download the 'tabular-data.zip' file.

There are two main options for accessing cover data:

(1) The seaviewsurvey_surveys.csv table contains percent cover of broad cover categories (hard coral, algae, soft coral, other invertebrate, and other), summarised among all photo-quadrats within a survey. This table has one row per survey (860 rows) and also contains the latitude and longitude of the start and end of the survey.

(2) There is regional variation in the cover classification schemes used (documented in the seaviewsurvey_labelsets.csv table). So if you need to access cover data at a more refined level of classification then the recommended procedure is:

Look at the seaviewsurvey_labelsets.csv table to understand how the classification scheme varies among regions, and develop a strategy for how you will need to group cover categories to best suite your application.

The five seaviewsurvey_reefcover_*.csv files contain the cover data for each of the five regions. These tables have one record per photo-quadrat, hence may contain large numbers of rows (the largest has over 300,000 rows). It is recommended that you use software that can handle tables with large number of rows to process these files, such as R, MySQL, PostGRES, Microsoft Access, etc. Manipulating these data with Microsoft Excel is not advised (especially the 32-bit version).

To create survey-scale summaries of the cover data you will need to: (i) summarise the reef cover table by "imageid", taking the mean cover value among all records for an image; (ii) summarise the resulting image summary data by "surveyid", taking the mean cover value among all records for a

survey. This two step process is recommended because some images have more than one $1m^2$ photo-quadrat image extracted from them, but these photo-quadrats are not independent. When merging cover classes (columns in these tables) you would need to sum the values from merged fields. This is a key point: always take the mean when combining rows, always sum when merging cover type columns. A good consistency check to use after merging classes is that the cover fields in each row should always sum to 1.0 (these values are proportions, and 1.0 represents 100%).

To georeference the survey summaries the recommended strategy is to join the resulting summary table to the seaviewsurvey_surveys.csv using they surveyid field, and access the latitude and longitude fields from that table. For more detailed spatial analysis the seaviewsurvey_reefcover_*.csv files contain the latitude and longitude for each photo-quadrat.


## Publications

Examples of peer-reviewed publications based on this dataset:

González-Rivero, M. *et al.* Scaling up Ecological Measurements of Coral Reefs Using Semi-Automated Field Image Collection and Analysis. *Remote Sensing* **8**, 30; (2016).

Peterson, E. E., E. Santos-Fernández, C. Chen, S. Clifford, J. Vercelloni, A. Pearse, R. Brown, B. Christensen, A. James, and K. Anthony. 2018. Monitoring through many eyes: Integrating scientific and crowd-sourced datasets to improve monitoring of the Great Barrier Reef. arXiv preprint arXiv:1808.05298.

Bryant, D. E. P., A. Rodriguez-Ramirez, S. Phinn, M. Gonzalez-Rivero, K. T. Brown, B. P. Neal, O. Hoegh-Guldberg, and S. Dove. 2017. Comparison of two photographic methodologies for collecting and analyzing the condition of coral reef ecosystems. Ecosphere **8**.

Beijbom, O., J. Hoffman, E. Yao, T. Darrell, A. Rodriguez-Ramirez, M. Gonzalez-Rivero, and O. H. Guldberg. 2015. Quantification in-the-wild: data-sets and baselines. arXiv preprint arXiv:1510.04811.

González-Rivero, M., P. Bongaerts, O. Beijbom, O. Pizarro, A. Friedman, A. Rodriguez-Ramirez, B. Upcroft, D. Laffoley, D. Kline, C. Bailhache, R. Vevers, and O. Hoegh-Guldberg. 2014. The Catlin Seaview Survey - kilometre-scale seascape assessment, and monitoring of coral reef ecosystems. Aquatic Conservation: Marine and Freshwater Ecosystems **24**:184-198.

Example of peer-reviewed papers relevant to the application of this dataset:

Beijbom, O. *et al.* Towards Automated Annotation of Benthic Survey Images: Variability of Human Experts and Operational Modes of Automation. *PloS one* **10**, e0130312; (2015).

Beijbom, O. *et al.* Improving automated annotation of benthic survey images using wide-band fluorescence. *Scientific reports* **6**, 23166; (2016).

Williams, Ivor Douglas, Courtney Couch, Oscar Beijbom, Thomas Oliver, Bernardo Vargas-Angel, Brett Schumacher, and Russell Brainard. "Leveraging automated image analysis tools to transform our capacity to assess status and trends on coral reefs." *Frontiers in Marine Science* 6 (2019): 222.

Vercelloni, J., S. Clifford, M. J. Caley, A. R. Pearse, R. Brown, A. James, B. Christensen, T. Bednarz, K. Anthony, M. González-Rivero, K. Mengersen, and E. E. Peterson. 2018. Using virtual reality to estimate aesthetic values of coral reefs. Royal Society Open Science **5**.

## License

These data are shared under the Creative Commons license "Attribution-NonCommercial-ShareAlike 4.0 International". A condition of the use of this data is that it is appropriately cited, attributed and derived datasets are shared under similar terms.

## Attribution for this dataset

When publishing using this dataset the following attribution should be included:

"Coral reef photo-quadrats and associated image classification data were developed by The University of Queensland using survey photographs sourced from the XL Catlin Seaview Suvey and XL Catlin Global Reef Record, undertaken by The University of Queensland and Underwater Earth, with additional funding from Vulcan Inc., Academica Sinica, and the Australian government."

**Attribution for the use of images in video or graphics contexts:**
"(c) Underwater Earth / XL Catlin Seaview Survey / The University of Queensland"

## Citation for this dataset

The citation for this dataset is:

González-Rivero, M., Rodriguez-Ramirez, A., Beijbom, O., Ganase, A., Dalton, P., Kennedy, E.V., Neal, B.P., Vercelloni, J., Bongaerts, P., Bryant, D.E.P., Brown, K., Kim, C., Radice, V.Z., Lopez-Marcano, S., Dove, S., Bailhache, C., Beyer, H.L. & Hoegh-Guldberg, O. 2019. Seaview Survey Photo-quadrat and Image Classification Dataset. The University of Queensland. DOI: 10.14264/uql.2019.930

## Acknowledgements

We acknowledge and thank the following funders of this work:

> Underwater Earth, with funding from XL Catlin (now AXA XL)
> Vulcan Inc.
> Australian Department for Energy and Environment
> Australian Research Council

We acknowledge and appreciate the contributions of the following people (alphabetical by surname):

| | |
|---|---|
| Norbert Englebert | Craig Hetherington |
| Patrick Gartrell | Tadzio Holtrop |
| Yeray Gonzalez-Marrero | Kathryn Markey |
| Susie Green | Angela Marulanda |
| Selene Gutierrez | Morana Mihaljevic |
| David Harris | Christopher Mooney |
| Kyra Hay | Sara Naylor |
| Ana Teresa Herrera-Reveles | Lorna Parry |

Coral bleaching at Lizard Island (Great Barrier Reef). The whitening of coral ("bleaching") occurs as a result of heat stress. © Underwater Earth / XL Catlin Seaview Survey / Christophe Bailhache